# MicroarrayTechnology

National Human Genome Research Institute

AFTER THE SEQUENCE:

WHOLE GENOME APPROACHES TO

BIOLOGICAL QUESTIONS

GENE EXPRESSION

GENE VARIATION

GENE FUNCTION

**PUBMED literature on DNA microarrays**

1050

**3.5 X average increase per year**

412

124

1    6    8    35

1995                                                    2001

*No. of publications per year*

---

National
Human
Genome
Research
Institute

Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |

Schena et al. Science 270:467

• Robotic high density printing of cDNAs

• Fluorescence detection

## Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |
|------|------|------|------|------|------|------|------|

DeRisi et al. Nat. Gen. 14:457

Schena et al. PNAS 14:1675

- Application to human cells

- Expression pattern related to tumorigenesis And T cell function

Lockhart et al. Nat. Biotech. 14:1675

- Oligonucleotides synthesized in situ

**Cancer Genetics Branch**

---

## Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |
|------|------|------|------|------|------|------|------|

Lakshari et al. PNAS 94:13057

Wodicka et al. Nat Biotech 13:1359

- Complete genome analysis: yeast

- Spotted DNA and oligos

**Cancer Genetics Branch**

## Development of Microarrays for Cancer Research

| 1995 | 1996 | 1997 | **1998** | 1999 | 2000 | 2001 | 2002 |
|------|------|------|------|------|------|------|------|

Khan et al. Cancer Res 58:5009

• Cancers of the same type cluster.

Eisen et al. PNAS 95:14863

• Two dimensional clustering.

Kononen et al. Nat Med 4:844
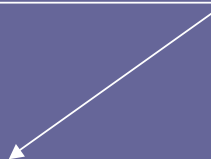
• Tissue microarrays.

Pinkel et al. Nat Gen 20:207

• CGH BAC arrays.

---

## Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | **1999** | 2000 | 2001 | 2002 |
|------|------|------|------|------|------|------|------|

Khan et al. PNAS 96:13464

• Expression program elicited by oncogene.

Golub et al. PNAS 286:531

• Formal diagnostic classifier.

Several sample clustering papers.

Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |

Bittner et al. Nature 406:536
- Class discovery within a cancer type.

Cancer Genetics Branch



Development of Microarrays for Cancer Research

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |

Bittner et al. Nature 406:536
- Class discovery within a cancer type.

Cancer Genetics Branch

Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |

Bittner et al. Nature 406:536
- Class discovery within a cancer type.

Alizadeh et al. Nature 406:503
- Class discovery correlating with outcome.

Perou et al. Nature 406:747
- Class discovery in breast cancer.

---

Development of Microarrays

| 1995 | 1996 | 1997 | 1998 | 1999 | 2000 | 2001 | 2002 |

Numerous publications addressing

- Class discovery and classification.
- Diagnostic classifiers.
- Biological/genetic correlations.
- Outcome correlations.
- Mathematical tools.

# MICROARRAY TERMINOLOG

- **Feature--an array element**

- **Probe--a feature corresponding to a defined sequence**

- **Target--a pool of nucleic acids of unknown sequence**

---

**Kinds of array elements**

- **Synthetic Oligonucleotides**

- **PCR products from**

  **Cloned DNAs**

  **Genomic DNA**

- **Cloned DNA**

# Microarray Manufacture

## • Printing

**Prepare cDNA probes**

"Normal"   Tumor

RT
Label with
Fluorescent Dyes

Combine
Equal
Amounts

Hybridize
probe to
microarray          SCAN

**Prepare microarray**

# Microarray Manufacture

- **Printing**
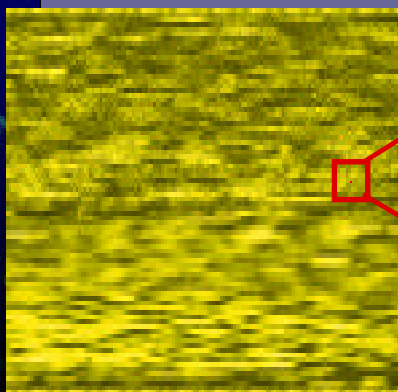
- **Synthesis *in situ***

# MICROARRAY READOUT

**Determine quantity of target bound to each probe in a complex hybridization**

- **Must have high sensitivity, low background**

- **High spatial resolution essential**

- **Dual channel capability preferred**

# DNA Microarray Applications

- **Resequencing**

  **Mutations**

  **Polymorphisms**

---

# *BRCA1* Coding Region Array

**V1713A**

T
G
C
A

5254 5255 5256 5257 5258 5259 5260

# Oligonucleotide Array Design

**Target**

C

**Surface Probes**

**Length**

A

25

C

25

G

25

T

25

| A |
|---|
| C |
| G |
| T |

**Perfect Match Probe**

---

# SINGLE NUCLEOTIDE POLYMORPHISM

## AGGTTACCAGTA

## AGGTTGCCAGTA

## OCCUR ABOUT 1: 1250 BASES

# SINGLE NUCLEOTIDE POLYMORPHISMS

- Polymorphic SNPs occur approximately every 1 kb in the human genome.

- Dense SNP maps provide a basis to design microarrays for genome scanning

---

# LABELLING SNPs

## Genomic DNA
↓ multiplex PCR

## Unlabeled amplicons
↓ primer extension

## Labeled amplicons
↓ pool, denature, dilute into buffer

## Hybridize to microarray

# SNP CHIP*



*Wang et al.
Science 280:1077
1998

A/A homozygote   A/C heterozygote   C/C homozygote

Query Base

Cutler DJ et al. Genome Res. 2001 11:1913-25.

# ACCURACY OF SNP CHIP

**Table 3.** ABACUS SNP Detection and Genotyping Accuracy

**A. Accuracy of autosomal SNPs detection**

|  | Verified | Total Possible |
|---|---|---|
| Singleton SNPs | 17 | 17 |
| Non-singleton SNPs | 91 | 91 |
| Total SNPs | 108 | 108 |

**B. Number of autosomal SNPs electronically verified**

| Number of SNPs electronically verified | 371 |
|---|---|

**C. Accuracy of autosomal genotype calls**

| Number of verified homozygous genotype calls | 1515 |
|---|---|
| Number of incorrect homozygous genotype calls | 0 |
| Percent correct homozygote calls | 100.00% |
| Number of verified heterozygous genotype calls | 423 |
| Number of incorrect heterozygous genotype calls | 3 |
| Percent correct heterozygote calls | 99.30% |

**D. Accuracy of haploid genotype calls**

| Number of bases sequenced (6X coverage) | 17,423 |
|---|---|
| Number of bases different from microarray chip calls | 0 |
| Percent of bases identical | 100.00% |

Cutler DJ et al. Genome Res. 2001 11:1913-25.

Cancer Genetics Branch

---

# SNP CHIP FOR ALLELIC IMBALANCE



Primdahl H et al. J Natl Cancer Inst. 2002, 94:216-223

Cancer Genetics Branch

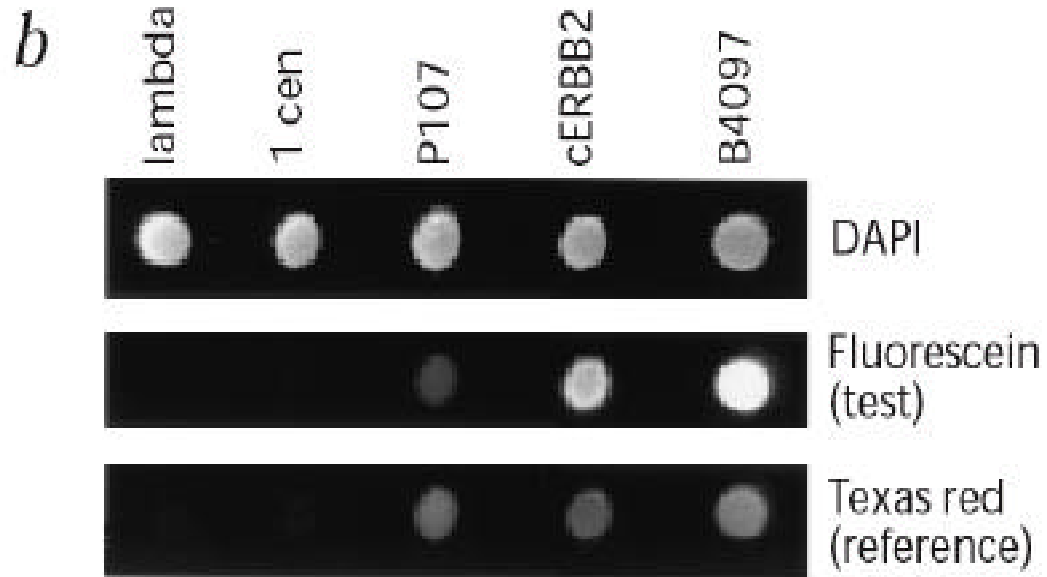# DNA Microarray Applications

- **Resequencing**

    **Mutations**
    **Polymorphisms**

- **Gene copy number**

---

# Gene Copy Number

- **Array format CGH**

- **Large insert clones**
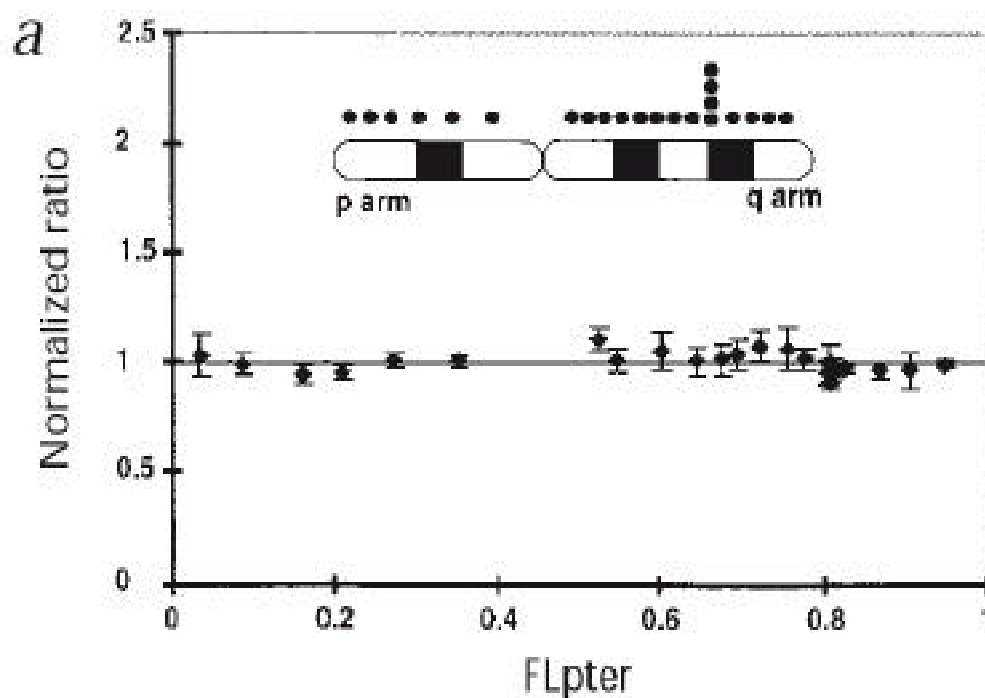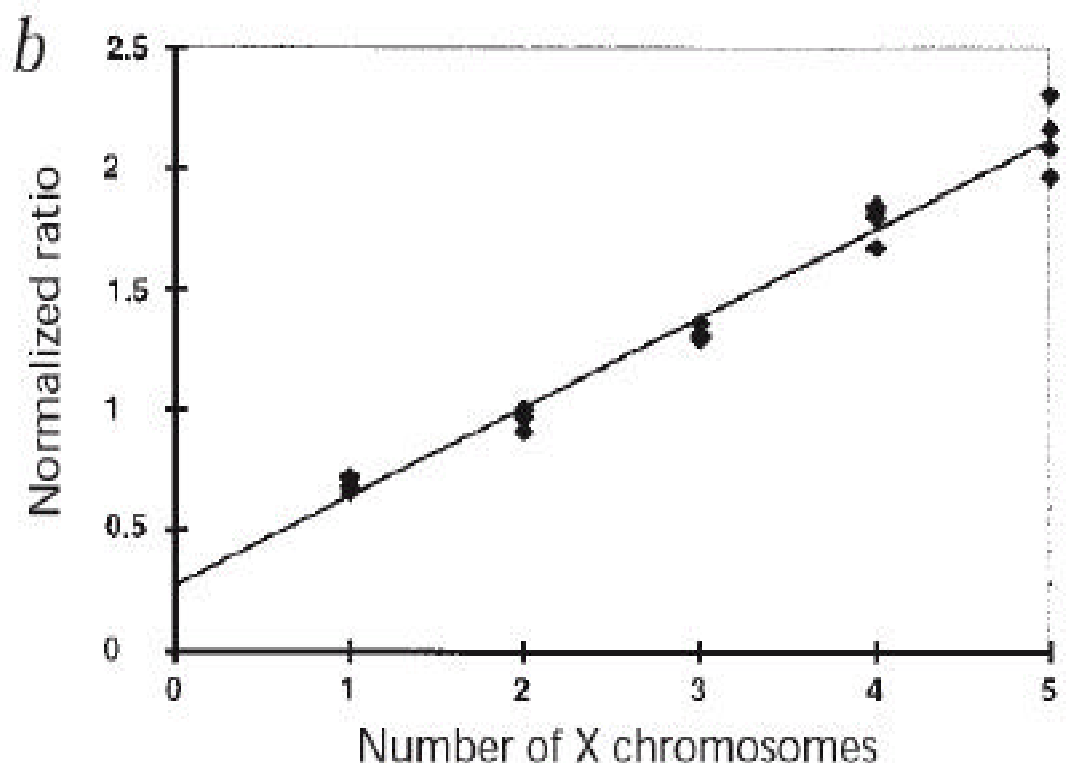
- **cDNA clones/exons**

# CGH BAC ARRAYS



Pinkel D et al., Nature Genetics 20, 207 - 211 ,1998.

Cancer Genetics Branch

# CGH BAC ARRAYS



Pinkel D et al., Nature Genetics 20, 207 - 211 ,1998.
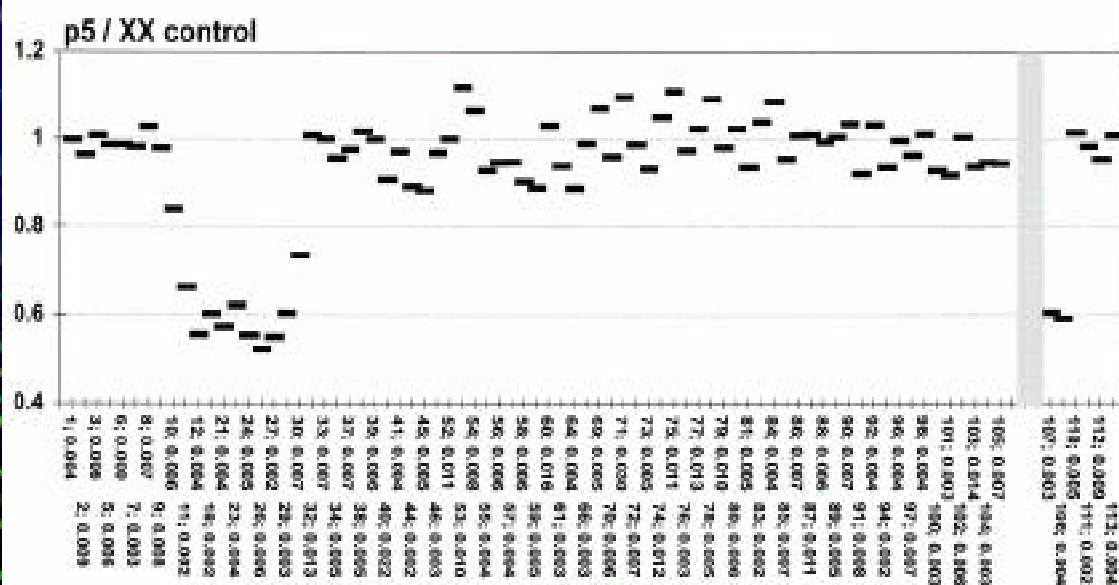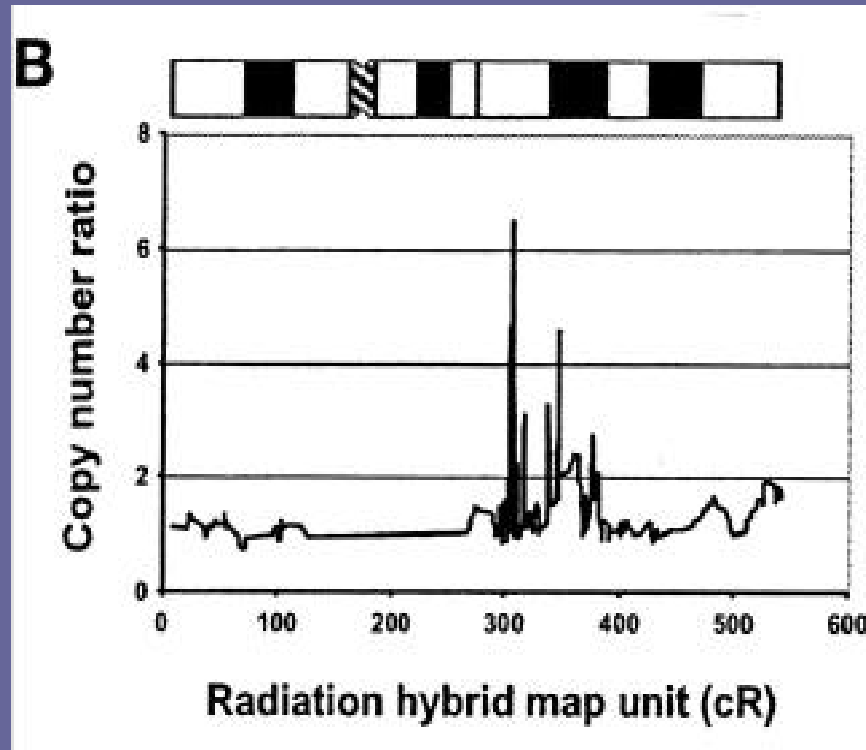
Cancer Genetics Branch

# CGH BAC ARRAYS



Pinkel D et al., Nature Genetics 20, 207 - 211 ,1998.

# CGH BAC ARRAYS



Bruder CE et al., Hum Mol Genet. 2001;10:271-82.

# CGH cDNA



Kauraniemi P et al., Cancer Res. 2001 ;61:8235-40.

# CGH cDNA



Kauraniemi P et al., Cancer Res. 2001

# DNA Microarray Applications

- **Resequencing**

  **Mutations**
  **Polymorphisms**

- **Gene copy number**

- **Gene expression**

---

# High throughput analysis of gene expression

- **cDNA library sequencing**

- **Serial analysis of gene expression (SAGE**

- **Microarray hybridization**

# STRATEGIES FOR SIGNAL GENERATION FROM mRNA

• Fluorochrome conjugated cDNA

• Ligand substituted nucleotides with  secondary dete

• Radioactivity

• RNA amplification

# Oligo versus cDNA Arrays for Expression Analysis

---

## Oligonucleotide Arrays: Pros

- Complete control over sequence

- Sequence and geometric perfection

- Extremely high feature density

# Oligonucleotide Arrays: Cons

- **Lack of Flexibility in Some Formats**

- **Absolute Requirement for Sequence Da**

- **Risk of uneven Performance by Individual Array Elements (Lack of Oligo Picking Rules)**

---

# cDNA Arrays: Pros

- **High Degree of Flexibility**

- **Sequence Independent**

- **High Stringency Hybridization**

- **High Signal Intensity: No Need for Signal Amplification**

# cDNA Arrays: Cons

- **Clone Availability**

- **Clone Handling**

- **Clone Authentication**

- **Possible Cross-hybridization**
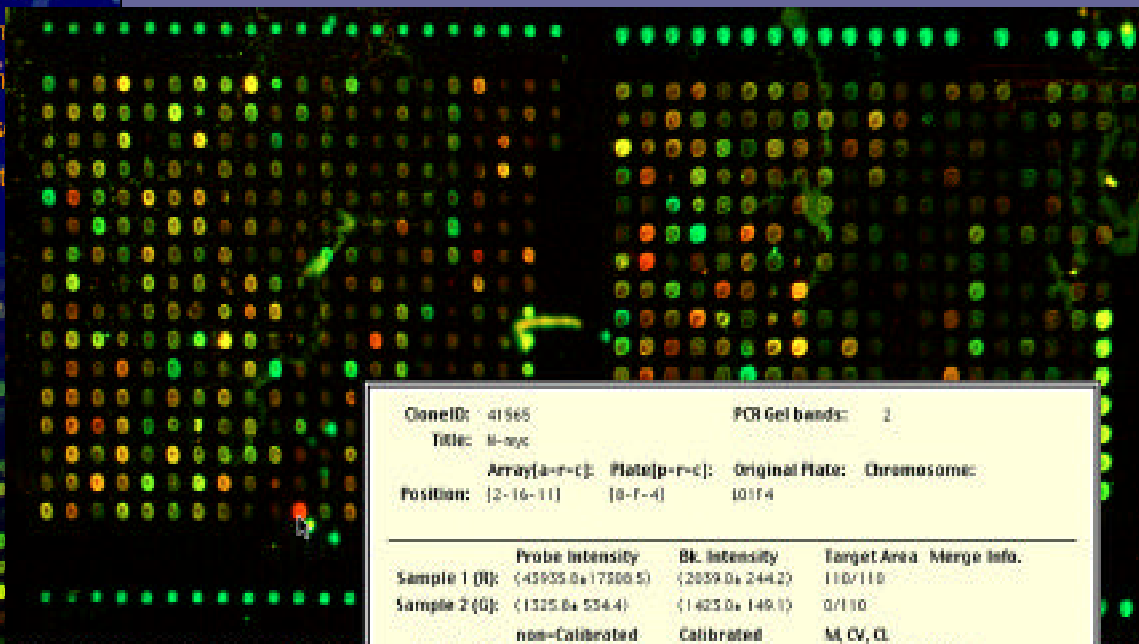
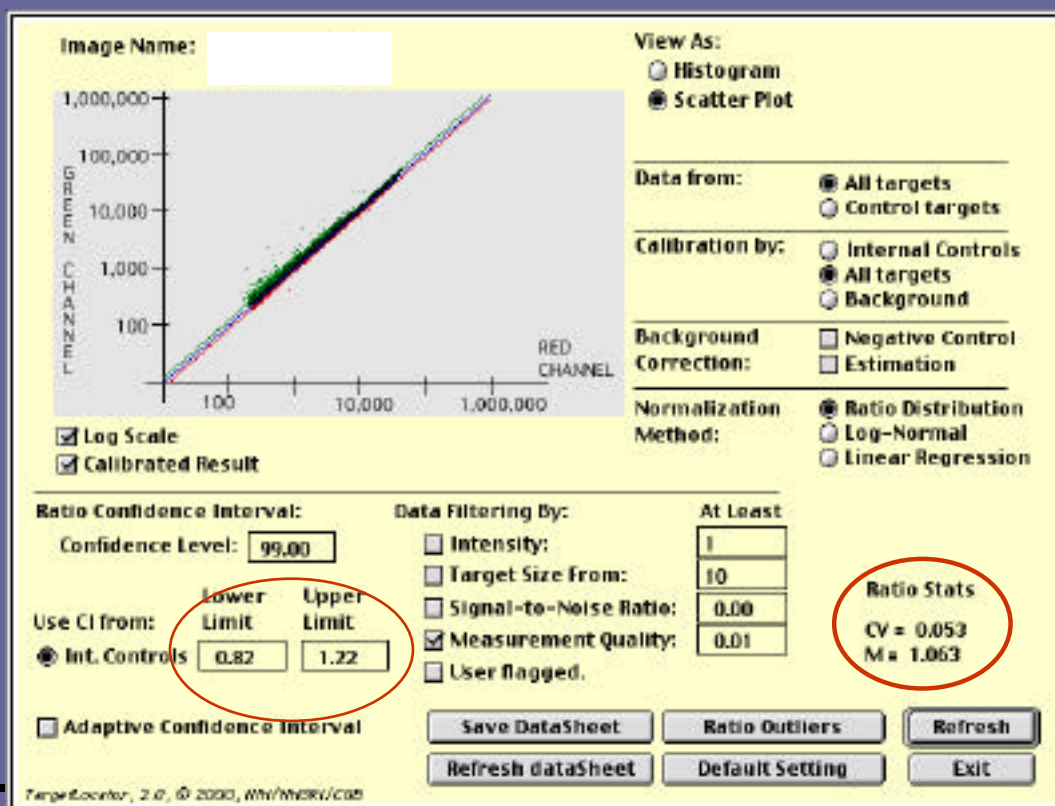# Image Analysis: DeArray

## Grid Overlay

## Target detection

# Image Analysis: DeArray

| | | | | |
|---|---|---|---|---|
| CloneID: | 41565 | | PCR Gel bands: | 2 |
| Title: | N-myc | | | |
| | Array[a-r-c]: | Plate[p-r-c]: | Original Plate: | Chromosome: |
| Position: | [2-16-11] | [0-F-4] | 101F4 | |

| | Probe Intensity | Bk. Intensity | Target Area | Merge Info. |
|---|---|---|---|---|
| Sample 1 (R): | (45935.0±17300.5) | (2039.0±244.2) | 110/110 | |
| Sample 2 (G): | (1325.0±554.4) | (1423.0±149.1) | 0/110 | |
| | non-Calibrated | Calibrated | M, CV, CL | |
| Ratio R/G: | 33.140 | 38.709 | 0.856, 0.258, 99.00% | |
| Interval: | [ 0.261, 2.825] | [ 0.305, 3.297] | | |

TargetLocator, 1.0, © 1997, NIH/NHGRI/CGB

## DATA QUALITY IS CRITICAL

---

**Output of cDNA microarrays: expression ratio**

**Output of oligonucleotide arrays:  expression lev**

**Both types of data can be analyzed with essentially the same tools.**

**Normalization of cDNA microarrays:**
•global
•housekeeping
•spiked standards

**Normalization of oligonucleotide arrays:**
•global

# Building Expression Arrays

- Completely sequenced and annotated organisms

- Complete gene list unavailable

# Expressed Sequence Tags

- Partial, inaccurate cDNA sequences

- Redundancy allows clustering by gene

- Incomplete representation of genome

# Clustering of Human ESTs

**Final Number of Clusters (sets) (Unigene 146)**

====================================

**96109    sets total**

**21857    sets contain at least one known gene**
**94916    sets contain at least one EST**
**20664    sets contain both genes and ESTs**

**• 1193 genes are not represented by an identifiable EST.**

---

## Histogram of cluster sizes for UniGene build 146

96,000 clusters

| Cluster size | Number of clusters |
|---|---|
| 1 | 33700 |
| 2 | 13467 |
| 3-4 | 15267 |
| 5-8 | 10197 |
| 9-16 | 5777 |
| 17-32 | 3894 |
| 33-64 | 3549 |
| 65-128 | 4031 |
| 129-256 | 3718 |
| 257-512 | 1732 |
| 513-1024 | 537 |
| 1025-2048 | 162 |
| 2049-4096 | 56 |
| 4097-8192 | 19 |
| 8193-16384 | 3 |

# Expressed Sequence Tags: Options for Array Constructi

- •"Standard" clone sets

- • Custom clone sets

- • Synthetic oligonucleotides

## Accessing Expression Data

- •Individual Lab and Journal Sites

# APPLICATIONS OF EXPRESSION ARRAYS

## •Direct comparisons (Induction)

### Biological system critical

## •Expression profiling

### Requires statistical tools

### Power arises from increasing sample number

---

# APPLICATIONS OF EXPRESSION ARRAYS

## •Statistical tools for large datasets

## •First generation approaches.

# APPLICATIONS OF EXPRESSION ARRAYS : TUMOR PROFILING

- Towards a molecular taxonomy of cancer

- Methods lead to gene identification

- Individualized diagnosis and therapy

# APPLICATIONS OF EXPRESSIONARRAYS : GENE IDENTIFICATION

- Groups of genes
  - Pathways
  - Co-regulated
  - Correlate with copy #
  - Correlate clinically

- Candidate disease genes

# APPLICATIONS OF EXPRESSION ARRAYS: TUMOR PROFILING

- Clustering
  - Unsupervised
  - Supervised
- Classification

- Can classify with respect to any clinically interesting variable

---

## Alveolar Rhabdomyosarcoma

**Pax3**
**chromosome 2**

# Method

- Compared 7 ARMS with 6 unrelated cancers cell lines

- Using cDNA microarray containing 1238 elements
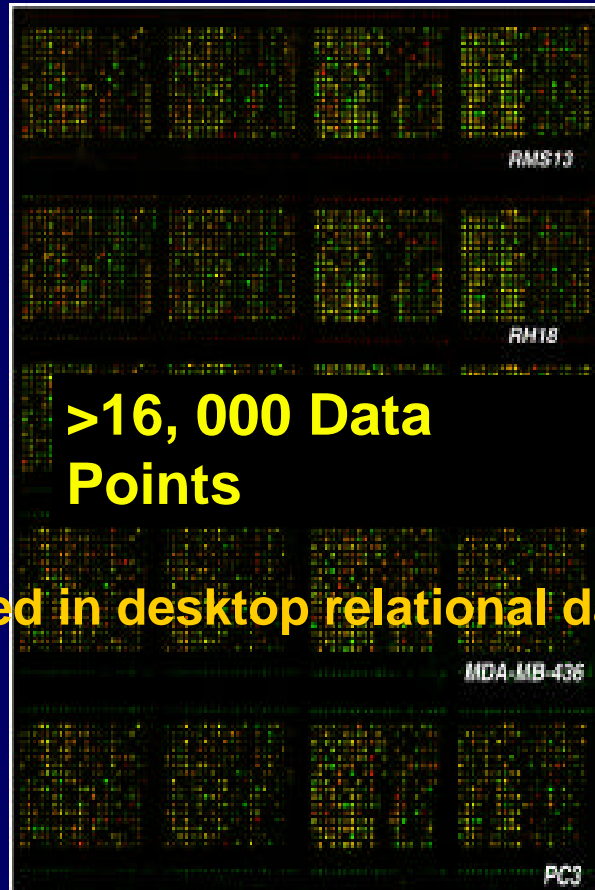
# Cell Line Characteristics

| Cell Line | Pax3-FKHR | Diagnosis |
|---|---|---|
| ARMS1 | + | ARMS |
| RH3 | + | ARMS |
| RH4 | + | ARMS |
| RH5 | + | ARMS |
| RMS18 | + | ARMS |
| RMS13 | + | ARMS |
| RH28 | + | ARMS |
| A204 | - | Undifferentiated Sarcoma |
| NGP127 | - | Neuroblastoma |
| TC71 | - | Ewing's Sarcoma |
| UACC-903 | - | Melanoma |
| PC3 | - | Prostate Carcinoma |
| MDA-MB-436 | - | Breast Carcinoma |
| Control | | |
| NIL-C | - | Fibroblast |

# Results
- **13 Experiments**
- **1238 genes**

*RMS13*

*RH18*

**>16, 000 Data Points**

**Data placed in desktop relational database**

*MDA-MB-436*

*PC3*

---

# SCATTER PLOT

## ARMS

| | RH4 | RH5 | RH3 | RH28 | RH18 | ARMS1 |
|---|---|---|---|---|---|---|
| **RMS13** | r=0.78 | **r=0.77** | r=0.69 | r=0.78 | r=0.73 | r=0.68 |
| **RMS13** | r=0.58 | **r=0.40** | r=0.48 | r=0.44 | r=0.46 | r=0.35 |
| | TC71 | PC3 | MDA-MB-435 | JACC-903 | NGP127 | A204 |

## NON-ARMS

# Matrix of Pearson Correlation Coefficients
## Distance Map
## 78 pair-wise comparisons

|  | RH3 | RH4 | RH5 | RMS13 | RH18 | **RH28** | A204 | NGP127 | TC71 | UACC-903 | MDA-MB-436 | PC3 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ARMS1 | 0.547 | 0.606 | 0.725 | 0.683 | 0.634 | 0.375 | 0.307 | 0.39 | 0.498 | 0.426 | 0.417 | 0.314 |
| **RH3** | | 0.760 | 0.726 | 0.69 | | **0.807** | | 0.565 | 0.566 | 0.391 | 0.452 | 0.403 |
| RH4 | | 0.771 | 0.778 | 0.67 | | 0.541 | 0.486 | 0.558 | 0.468 | 0.555 | 0.476 |
| RH5 | | | 0.769 | 0.667 | 0.751 | 0.37 | 0.486 | 0.607 | 0.43 | 0.532 | 0.447 |
| RMS13 | | | | 0.731 | 0.746 | 0.35 | 0.463 | 0.582 | 0.446 | 0.475 | 0.404 |
| RH18 | | | | | 0.703 | 0.274 | 0.281 | 0.549 | 0.389 | 0.405 | 0.36 |
| RH28 | | | | | | 0.417 | 0.493 | 0.644 | 0.479 | 0.478 | 0.42 |
| A204 | | | | | | | 0.426 | 0.361 | 0.398 | 0.368 | 0.377 |
| NGP127 | | | | | | | | 0.352 | 0.241 | 0.371 | 0.368 |
| TC71 | | | | | | | | | 0.46 | 0.456 | 0.472 |
| UACC-903 | | | | | | | | | | 0.507 | 0.538 |
| MDA-MB-436 | | | | | | | | | | | 0.662 |
| PC3 | | | | | | | | | | | |

## Hierarchical Clustering Dendrogram

G
E
N
E
S

Cancer Genetics Branch



Cancer Genetics Branch

# Multidimensional Scaling

|          | Seattle | L.A. | Houston | D.C.  |
|----------|---------|------|---------|-------|
| Seattle  | 0       |      | 1,200   | 2,500 |
| L.A.     |         | 0    | 1,500   | 2,800 / 2,600 |
| Houston  |         |      | 0       | 1400  |
| D.C.     |         |      |         | 0     |

|      | Exp1 | Exp2 | Exp3 | Exp4 |
|------|------|------|------|------|
| Exp1 | 0    | 0.50 | 0.25 | 0.05 |
| Exp2 |      | 0    | 0.50 | 0.55 |
| Exp3 |      |      | 0    | 0.60 |
| Exp4 |      |      |      | 0    |

Seattle
D.C.
L.A.
Houston

Exp4
Exp1
Exp2
Exp3

---

## Multidimensional Scaling Analysis

NGP127
A204
RH3
RH28
RH4
ARMS1
MDA-MB436
RH5   RMS13
PC3
RH18
TC71
UACC-903

CLASS DISCOVERY IN MELANOMA

Cancer Genetics Branch



How subtle a difference in pattern could be reliably detected?

The similarity of the results of the 8 duplicate experiments would put all of these samples within the volume of the central blue sphere.

Cancer Genetics Branch

Weighted List

ARTICLES

# Classification and diagnostic prediction of cancers using gene expression profiling and artificial neural networks

JAVED KHAN[1,2], JUN S. WEI[1], MARKUS RINGNÉR[1,3], LAO H. SAAL[1], MARC LADANYI[4], FRANK WESTERMANN[5], FRANK BERTHOLD[6], MANFRED SCHWAB[5], CRISTINA R. ANTONESCU[4], CARSTEN PETERSON[3] & PAUL S. MELTZER[1]

[1]Cancer Genetics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland, USA
[2]Pediatric Oncology Branch, Advanced Technology Center, National Cancer Institute, Gaithersburg, Maryland, USA
[3]Complex Systems Division, Department of Theoretical Physics, Lund University, Lund, Sweden
[4]Department of Pathology, Memorial Sloan-Kettering Cancer Center, New York, New York, USA
[5]Department of Cytogenetics, German Cancer Research Center, Heidelberg, Germany
[6]Department of Pediatrics, Klinik für Kinderheilkunde der Universität zu Köln, Köln, Germany
J.K., J.S.W. and M.R. contributed equally to this study.
Correspondence should be addressed to J.K. or P.S.M.; email: khanjav@mail.nih.gov or pmeltzer@nhgri.nih.gov

# Molecular Taxonomy of Small Round Blue Cell Tumors

## Hypothesis

- Using cDNA microarrays we can identify
  the genes whose expression level is
- Utilize these genes to classify the characteristic for that cancer & type
  small blue round cell tumors into the
  correct diagnostic categories

# Model: Small Blue Round Cell Tumors

**Ewing's**

**Neuroblastoma**

**Rhabdomyosarcoma**

**Lymphoma**

# cDNA Microarray Analysis

## Experiments

| | |
|---|---|
| Burkitt's Lymphoma | 8 |
| EWS-Tumor | 13 |
| EWS-Cell line | 10 |
| Neuroblastoma | 12 |
| RMS-Tumor | 10 |
| RMS-Cell line | 10 |
| Test Unknown | 25 |

**>500, 000 Data Points**

**cDNA microarray**
**4,000 sequence verified known genes**
**2,567 sequence verified EST**
**6,567 genes**

---

Artificial Neural Networks

Pattern Recognition

Training

## ANN Output
## (96 Genes)

**EWS**

Classification 100%

**RM**
**S**

**NB**

**BL**

Table 1 Training sample characteristics

---

# Sensitivity Measurement
# Ranked Genes

**High Sensitivity**
**High Rank**

1

0

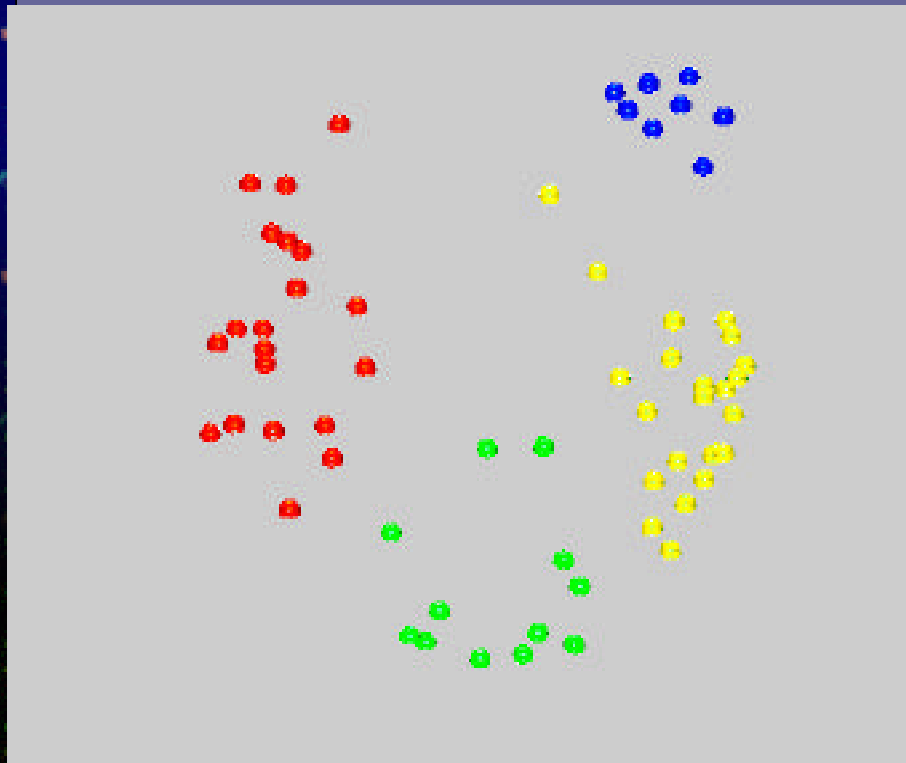ANN Output

Gene Expression

# Sensitivity Measurement Ranked Genes



# Gene Minimization

## MULTIDIMENSIONAL SCALING

- Lymphoma
- RMS
- NBL
- EWS

Cancer Genetics Branch



## Artificial Neural Networks

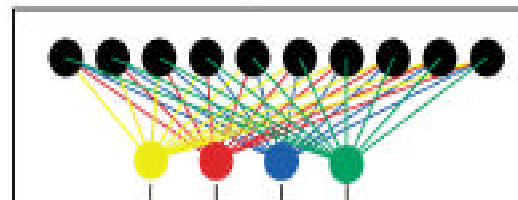## Recalibrate with Top 96 Genes

## Diagnosis?

63 training+ 25 "unknown"

96 genes

↓

PCA (10)

↓

Input

EWS  RMS  BL  NB

Output

# Diagnostic Classification



---

WHAT HAVE WE LEARNED FROM THE EXPRESSION PROFILING OF CANCERS SO FAR?

- DISTINCT HISTOLOGIES HAVE DISTINCT PATTERNS OF GENE EXPRESSION.

- USING EXPRESSION DATA IT IS POSSIBLE TO DEVELOP ROBUST FORMAL DIAGNOSTIC CLASSIFIERS.

- NOVEL SUBGROUPS CAN BE RECOGNIZED WITHOUT PREVIOUSLY DEFINED HISTOPATHOLOGIC CORRELATES.

- CLINICALLY USEFUL CATEGORIES CAN BE DEFINED, BUT DO NOT NECESSARILY REQUIRE ARRAYS FOR EVERYDAY CLINICAL IDENTIFICATION.

National
Human
Genome
Research
Institute

WHAT WE HOPE TO LEARN IN THE FUTURE

• IMPROVE THE DIAGNOSTIC CATEGORIZATION
OF TUMORS.

• IDENTIFY USEFUL PREDICTIVE MARKERS FOR
OUTCOME AND THERAPEUTIC RESPONSE
(ARRAY OR CONVENTIONAL).

• IDENTIFY POINTS FOR INTERVENTION:

CRITICAL PATHWAYS

DRUG TARGETS

---

CLINICAL CORRELATIVE STUDIES
USING MICROARRAYS

• DEFINE QUESTION AND PATIENT SAMPLE.

• APPROPRIATE AND RIGOROUS STATISTICAL
ANALYSIS OF ARRAY DATA.

• RESULT: GENES WHICH CARRY INFORMATION
RELEVANT TO QUESTION POSED.

• DEVELOP FORMAL CLASSIFIER.

• VALIDATE ON ADDITIONAL SAMPLE SET.

## MODEL SYSTEM WITH CLEAR THERAPEUTIC IMPLICATIONS:
## GASTROINTESTINAL STROMAL TUMOR

- RELATED TO THE INTERSTITIAL CELLS OF CAJAL

- KIT MUTATIONS

- STI-571 SENSITIVITY

- THE BEST "CREDENTIALED" TARGETS ARE THOSE ACTIVATED BY MUTATION.
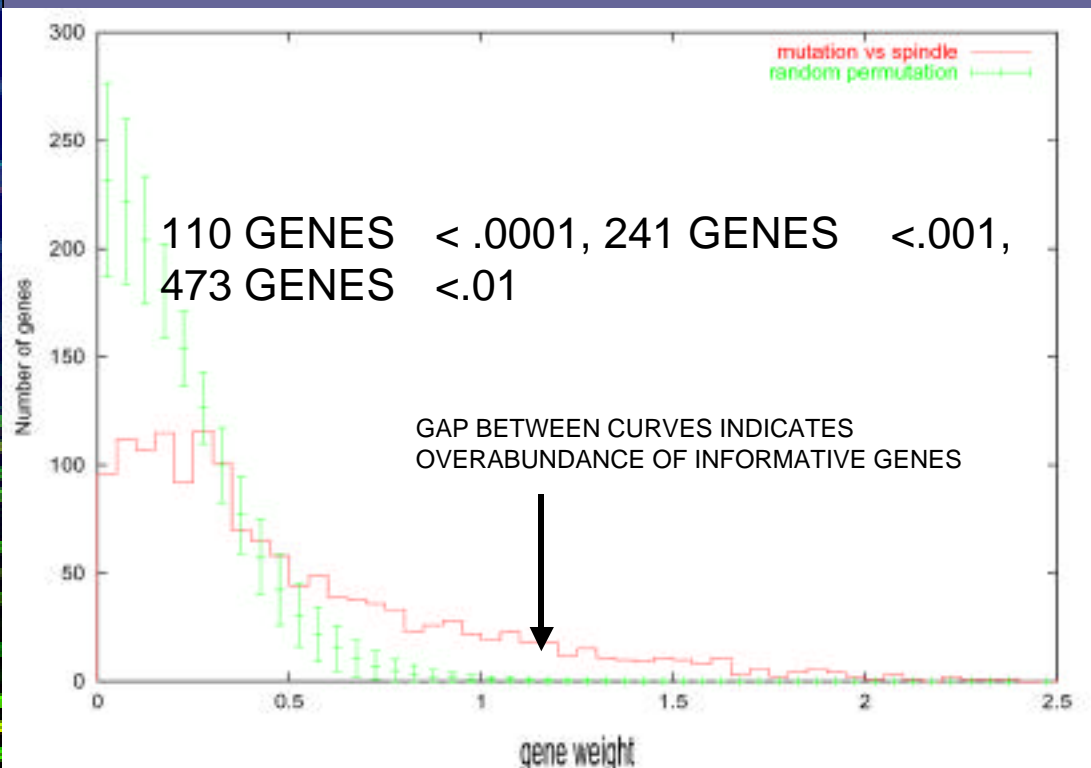
## SAMPLES

- 13 MALIGNANT GISTs

- ALL KIT POSITIVE BY IHC

- ALL WITH KIT MUTATIONS

- 4 GUT WALL PRIMARIES

- 8 INTRA-ABDOMINAL EXTENSION

- 1 LIVER METASTASIS

- 6 COMPARISON TUMORS: EXTRA-GI SPINDLE CELL MORPHOLOGY

## ARRAYS AND DATA ANALYSIS

- 13,824 ELEMENT SPOTTED cDNA ARRAYS

- OSA REFERENCE PROBE

- GENES RANKED FOR EXPRESSION IN GIST

- WEIGHTED DISCRIMINATOR LIST

- MDS AND HIERARCHICAL CLUSTERING

## OVERABUNDANCE OF INFORMATIVE GENES



110 GENES  < .0001, 241 GENES  <.001, 473 GENES  <.01

GAP BETWEEN CURVES INDICATES OVERABUNDANCE OF INFORMATIVE GENES
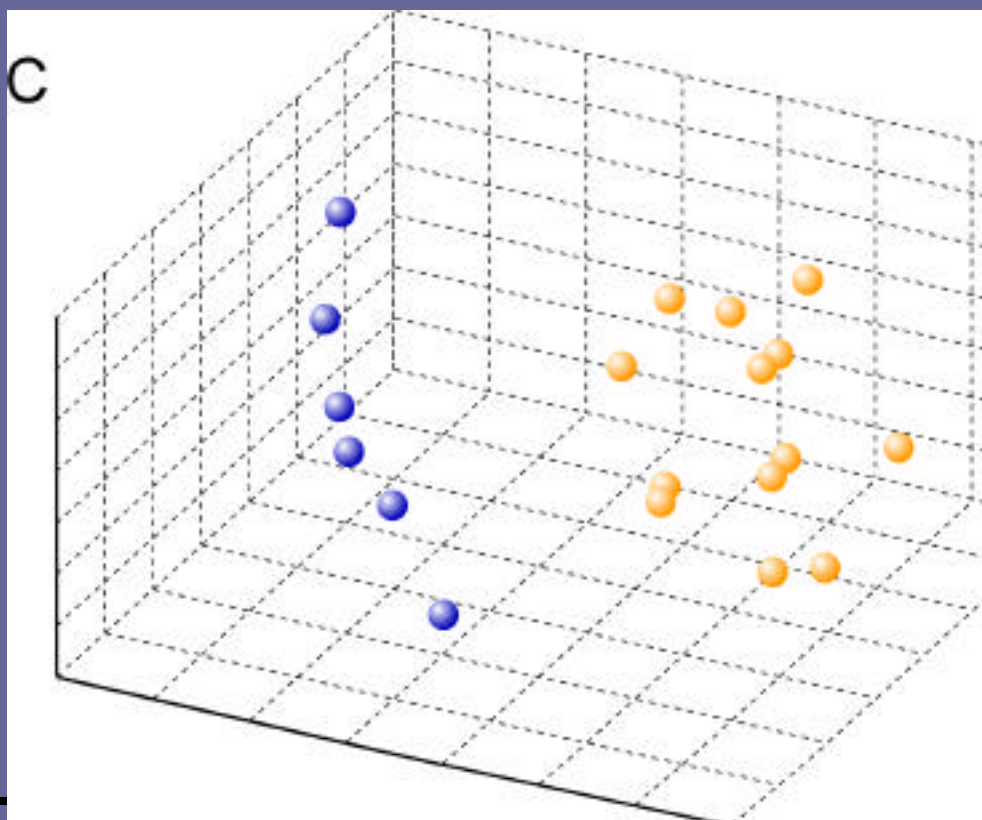
## DATA ANALYSIS

- PREFILTER FOR QUALITY AND IDENTIFY GENES HIGHLY EXPRESSED IN GIST: 1987 GENES

- RANK BY WEIGHTED DISCRIMINATOR METHOD
  $w(g, \pm) = \mu_{+}(g) - \mu_{-}(g) \, /[\sigma_{+}(g) + \sigma_{-}(g)]$

- RANDOM PERMUTATION TEST ($10^5$ trials)
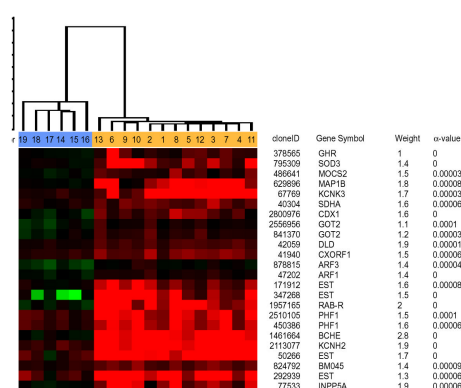
- CLUSTER ANALYSIS USING DISCRIMINATORS

---

## MDS PLOT



C

# TOP DISCRIMINATORS FOR GIST

| Rank | Weight | Alpha | Gene Description |
|------|--------|-------|------------------|
| 1 | 7.55575 | 0 | v-kit sarcoma oncogene |
| 2 | 6.48306 | 0 | G coupled receptor 20 |
| 3 | 4.60057 | 0 | G coupled receptor 20 |
| 4 | 4.51681 | 0 | annexin A3 |
| 5 | 3.33057 | 0 | KIAA0353 protein |
| 6 | 3.31734 | 0 | phosphofructokinase, muscle |
| 7 | 2.95095 | 0.00008 | DKFZP434N161 protein |
| 8 | 2.83435 | 0 | protein kinase C, theta |
| 9 | 2.79721 | 0 | butyrylcholinesterase |
| 10 | 2.72752 | 0 | annexin A3 |

| | Gene Symbol | Weight | α-value |
|---|---|---|---|
| | PFKM | 3.3 | 0 |
| 489626 | DMN | 3.3 | 0 |
| 1161564 | ANXA3 | 4.5 | 0 |
| 2469213 | GPR20 | 6.5 | 0 |
| 2568905 | KIT | 7.6 | 0 |
| 269806 | PTP4A3 | 2.6 | 0 |
| 375827 | SCG2 | 2.6 | 0 |
| 174627 | PRKCQ | 2.8 | 0 |
| 205239 | TNFRSF6B | 2 | 0 |
| 2832322 | PRKCQ | 2.6 | 0 |
| 2164126 | | | |

## CONCLUSIONS

• MALIGNANT GISTS EXHIBIT A DISTINCT AND HIGHLY COHERENT GENE EXPRESSION PROFILE.

• THIS GENE LIST IS RELEVANT BOTH TO GIST GROWTH AND NORMAL ICC FUNCTION.

• KIT, A CRITICAL GENE IN REGULATING GIST GROWTH, IS THE BEST DISCRIMINATOR FOR THIS DISEASE.

• EXTENDING THIS APPROACH TO OTHER CANCERS MAY HELP IDENTIFY NEW DISEASE SPECIFIC DRUG TARGETS.

---

# A RECURRING PROBLEM

**Oncogenes**

**Transcription factors**

**Hormones/growth factors**

**Drugs**

**Toxins**

**Radiation**

**?????**

**Downstream Genes**

•**Direct targets**

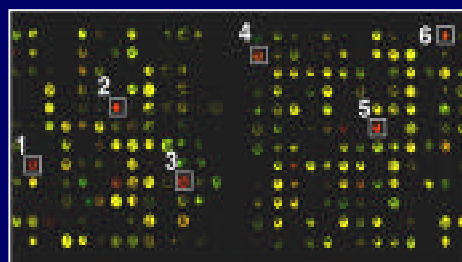•**Indirect targets**

**Retrovirus es:**

**Empty vector**
**Pax3**
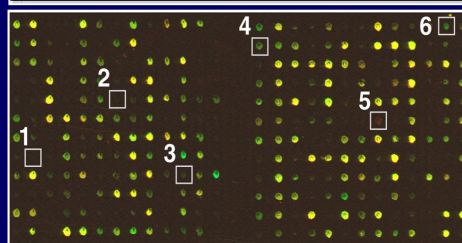**Pax3-Fkhr**

---

# RESULTS

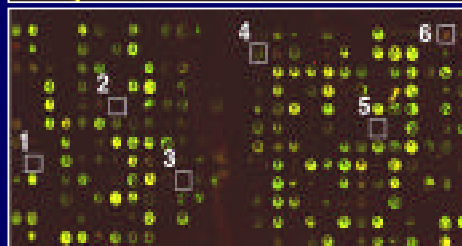**Mouse cDNA array 2200 genes**

**PAX3-FKHR**
**vs**
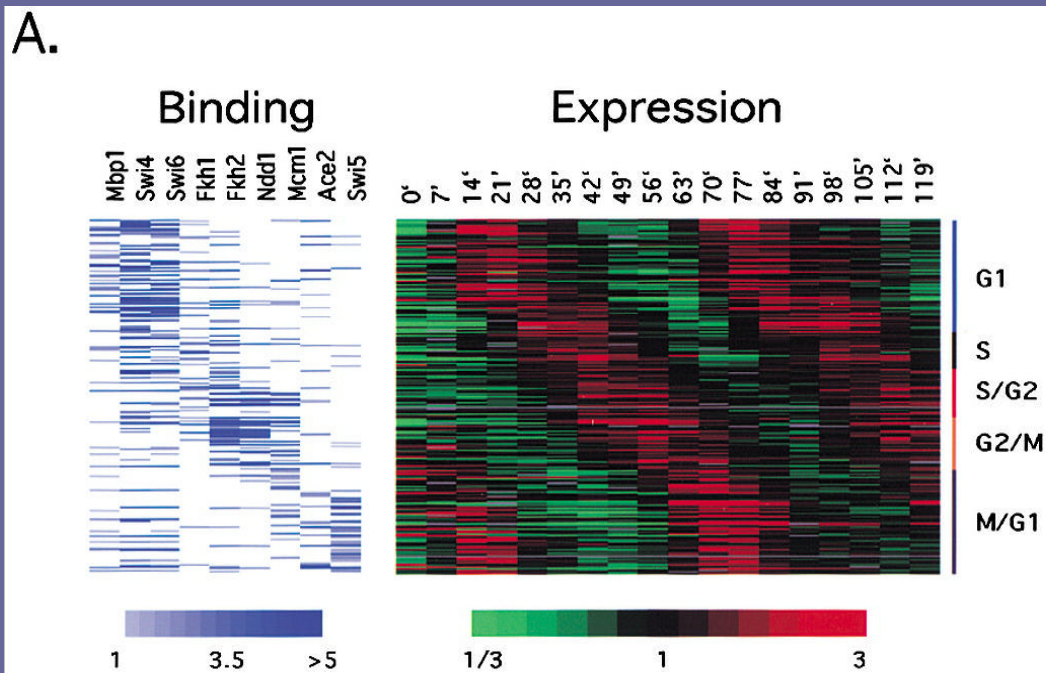**Empty Vector**



**PAX3**
**vs**
**Empty Vector**



**3T3 PARENT**
**vs**
**Empty Vector**



1. Troponin C
2. IGFBP5
3. Myogenin
4. Six1
5. Troponin T
6. IGF2

## Promoter Occupancy During Yeast Cell Cycle



Simon I Cell. 2001 Sep 21;106(6):697-708.

---

# Future Promise

• **Numerous clever applications to additional systems, new
computational tools and technical innovations.**
• **Development of large clinically annotated datasets
which allow precise definition of patient subsets and
lead to the identification of new therapeutic targets.**

• **Linking genomic sequence to expression data to define
regulatory elements/transcription factors associated with
co-regulated genes.**

•**Development of computational tools which allow
predictive
modeling of gene networks.**
•**Introduction of technologies which increases the
dimensionality of expression data: protein arrays; cell
arrays; promoter arrays.**